

Necessity of Numerical Smoothness

Tong Sun

Department of Mathematics and Statistics

Bowling Green State University

Bowling Green, OH 43403

Abstract. Numerical solutions of differential equations are usually not smooth functions. However, they should resemble the smoothness of the corresponding real solutions in one way or another. In [2] and [3], a kind of spacial smoothness indicators was defined and subsequently applied on the *a posteriori* error analysis. Here we prove that the boundedness of those smoothness indicators is actually a necessary condition for a piecewise polynomial function to approximate a smooth function with optimal convergence rate. This should help in validating the error analysis in [2] and [3]. Moreover, the result of this paper provides an efficient practical method to detect the loss of convergence rate due to the lack of numerical smoothness, hence it serves as a criterion for the qualities of many numerical schemes.

keywords. Numerical smoothness, smoothness indicator, necessary condition.

AMS subject class. Primary: 65M12 Secondary: 65M15

1 Introduction

In the literature of numerical solutions of partial differential equations, the smoothness of numerical solutions has not been a popular concept. Of course, numerical solutions obtained by using finite difference, finite element and finite volume methods are typically not smooth functions. Therefore, it seems on the surface that there is no smoothness whatsoever. However, it is widely known that a scheme for solving a time dependent problem needs to be dissipative and/or total variation diminishing. As a matter of fact, these are actually concerning the smoothness of numerical solutions. Although nonsmooth solutions (shocks, interfaces, etc.) are sometimes of great interest, PDE solutions are usually at least piecewise smooth. A numerical scheme must be able to approximate smooth pieces of solutions. Being dissipative and total variation diminishing is certainly necessary. However, as shown in this paper, there are actually stronger necessary conditions to be satisfied if a scheme is hoped to converge at its desired optimal rate.

In [2] and [3], a kind of spatial smoothness indicators is proposed. These indicators are subsequently verified to be bounded in the numerical experiments, which means that the computed numerical solutions are “numerically smooth”. Most importantly, these smoothness indicators are applied to the local error analysis, playing the role of higher order derivatives. For the scalar nonlinear conservation laws studied in [2] and [3], obtaining the local error estimates in terms of the numerical smoothness made it possible to do error propagation analysis by directly using the L^1 -contraction between entropy solutions. Consequently, *a posteriori* error estimates of optimal convergence rates and linear growth were obtained for the RK-DG scheme and the WENO scheme. The key advantage of the methodology of [2] and [3] is that the smoothness indicators serve as a bridge for bypassing the difficulty of proving any global property of a scheme, caused by the nonlinearities.

Here, we present the spatial smoothness indicator used in [3]. We prove that the boundedness of the smoothness indicator is actually a necessary condition for the numerical solution to converge to any smooth function, including the real solution, at the optimal convergence rate. In fact, the boundedness of the smoothness indicators is our choice of a global property to deal with in [2] and [3]. Since the numerical solutions are computed by the very complex DG, WENO and Runge-Kutta schemes, we are certainly unable to give any *a priori* proof for the boundedness of the smoothness indicators. The computation of the smoothness indicators is actually the way to bypass the difficult proof. The necessity result of this article confirms that it is reasonable to expect the boundedness of the smoothness indicators. The result not only further supports the error analysis of [2] and [3], but also supports the general methodology developed over there.

Beyond validating the methodology of [2] and [3], the necessity result provides an extremely efficient criterion on numerical solutions of differential equations. Namely, if the smoothness indicator computed from a numerical solution is too large, then the numerical solution is certainly not converging to any smooth function at the desired optimal rate. In another word, whenever the smoothness indicators seem too large, either there is some kind of non-smoothness in the PDE solution being approximated, or there is something wrong in the numerical scheme being used.

The result is formulated in a 1-D uniform partition for the simplicity. One can generalize the proof to higher dimensions and non-uniform triangulations. However, as the first result of its kind, the investigation of a simple 1-D result suffices to show the necessity of numerical smoothness.

2 The main results

Let \mathcal{P} be the space of polynomials of degree p or less, and \mathcal{P}_h be the space of piecewise polynomials of degree p or less on the uniform partition $a = x_0 < x_1 < \dots < x_N = b$, $h = x_{i+1} - x_i = (b - a)/N$. Let u^R be a piecewise polynomial in \mathcal{P}_h . The reason for using the superscript R is to indicate that u^R might be a reconstructed numerical solution, possibly from nodal values, cell averages, or other forms of numerical solutions. Define a smoothness indicator for any $u^R \in \mathcal{P}_h$ as

$$S^p = (\bar{M}, \bar{D})$$

with

$$\bar{M} = (\tilde{M}_0, \tilde{M}_1, \dots, \tilde{M}_{N-1}), \quad \text{where} \quad \tilde{M}_i = (M_i^0, M_i^1, \dots, M_i^p),$$

and

$$\bar{D} = (\tilde{D}_1, \dots, \tilde{D}_{N-1}), \quad \text{where} \quad \tilde{D}_i = (D_i^0, D_i^1, \dots, D_i^p).$$

Furthermore,

$$M_i^k = \frac{d^k}{dx^k} u^R(x_i^+), \quad L_i^k = \frac{d^k}{dx^k} u^R(x_i^-), \quad J_i^k = M_i^k - L_i^k$$

and, as in [3],

$$D_i^k = J_i^k / h^{p+1-k}.$$

In this article, W_q^k is the standard notation of the Sobolev space of the functions, where the k -th derivative is L^q -integrable. $H^k = W_2^k$. See [1], Chapter 2.

Definition 2.1 *A piecewise polynomial $u^R \in \mathcal{P}_h$ is numerically W_∞^{p+1} -smooth in the partition of cell size h , if there is a constant M_∞ such that all the components of S^p are bounded by M_∞ ; u^R*

is numerically H^{p+1} -smooth in the partition if there is a constant M_2 such that $\sum_{i=1}^{N-1} h [(D_i^0)^2 + (D_i^1)^2 + \dots + (D_i^p)^2] \leq M_2^2$; u^R is numerically W_1^{p+1} -smooth in the partition if there is a constant M_1 such that $\sum_{i=1}^{N-1} h [|D_i^0| + |D_i^1| + \dots + |D_i^p|] \leq M_1$.

It is obvious that, for the piecewise polynomial u^R to approximate a smooth function in the interior of all the cells in the domain, one expects $|M_i^k| \leq \mathcal{O}(1)$ for all k and i . It is also obvious that, if u^R is an approximation of a smooth function with the optimal convergence rate, one can expect $|D_i^0| \leq \mathcal{O}(1)$ for all i . In the case $k = p$, it is obvious that $|D_i^p| \leq \mathcal{O}(1)$ implies that the piecewise constant function $\frac{d^p}{dx^p} u^R(x)$ has bounded variation. In the main theorem below, we show the necessity of the boundedness of the other components of the smoothness indicator. Let's start with a simple lemma.

Lemma 2.2

$$Q(D^0, D^1, \dots, D^p) = \min_{\hat{v} \in \mathcal{P}} \left(\left\| \hat{v} + \frac{1}{2} \sum_{k=0}^p \frac{D^k}{k!} \xi^k \right\|_{L^2(-\frac{1}{2}, 0)}^2 + \left\| \hat{v} - \frac{1}{2} \sum_{k=0}^p \frac{D^k}{k!} \xi^k \right\|_{L^2(0, \frac{1}{2})}^2 \right)$$

is a positive definite quadratic form.

Proof. In the process of calculating the coefficients of the polynomial $\hat{v} = \sum_{k=0}^p V_k \xi^k$ to attain the minimum, each V_k must be a linear combination of D^0, D^1, \dots, D^p , hence the minimum is a quadratic form of D^0, D^1, \dots, D^p . The quadratic form is positive definite because, if any entry D^k of (D^0, D^1, \dots, D^p) is non-zero, then we have either $V_k + \frac{1}{2} \frac{D^k}{k!} \neq 0$ or $V_k - \frac{1}{2} \frac{D^k}{k!} \neq 0$. Due to the linear independence of $\{1, \xi, \dots, \xi^p\}$, one of the two terms in the minimum must be positive. #

Theorem 2.3 Suppose that $u \in H^{p+1}(a, b)$, and u^R is a piecewise polynomial function described as above. Then, there is a constant $C_2 > 0$, independent of h , u and u^R , such that

$$\|u - u^R\|_{L^2(a, b)} \geq h^{p+1} \left[\sqrt{\sum_{i=1}^{N-1} h Q(D_i^0, D_i^1, \dots, D_i^p)} - C_2 |u|_{H^{p+1}(a, b)} \right]. \quad (2.1)$$

Proof. Let \mathcal{P}_c be the space of piecewise polynomials of degree p or less on the partition $a < x_{\frac{1}{2}} < x_{\frac{3}{2}} < \dots < x_{N-\frac{1}{2}} < b$, where $x_{i-\frac{1}{2}} = (x_{i-1} + x_i)/2$. Since $u \in H^{p+1}(a, b)$, there is a constant $C_2 > 0$, independent of h and u , and a piecewise polynomial $u^I \in \mathcal{P}_c$, such that

$$\|u - u^I\|_{L^2(a, b)} \leq C_2 h^{p+1} |u|_{H^{p+1}(a, b)}. \quad (2.2)$$

Now, since $\|u^I - u^R\|_{L^2(a, b)} \leq \|u^I - u\|_{L^2(a, b)} + \|u - u^R\|_{L^2(a, b)}$, we have

$$\begin{aligned} \|u - u^R\|_{L^2(a, b)} &\geq \|u^I - u^R\|_{L^2(a, b)} - \|u - u^I\|_{L^2(a, b)} \\ &\geq \sqrt{\sum_{i=1}^{N-1} \|u^I - u^R\|_{L^2(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})}^2} - C_2 h^{p+1} |u|_{H^{p+1}(a, b)}. \end{aligned} \quad (2.3)$$

Let $Q_i^k = (M_i^k + L_i^k)/2$ and $w(x) = \sum_{k=0}^p \frac{Q_i^k}{k!} (x - x_i)^k$. Then, $u^R - w = \frac{1}{2} \sum_{k=0}^p \frac{J_i^k}{k!} (x - x_i)^k$ for $x \in \Delta_i^+ = (x_i, x_{i+\frac{1}{2}})$, and $u^R - w = -\frac{1}{2} \sum_{k=0}^p \frac{J_i^k}{k!} (x - x_i)^k$ for $x \in \Delta_i^- = (x_{i-\frac{1}{2}}, x_i)$.

$$\begin{aligned}
& \|u^I - u^R\|_{L^2(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})}^2 \geq \min_{v \in \mathcal{P}} \|v - u^R\|_{L^2(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})}^2 \\
& = \min_{v \in \mathcal{P}} \|v - (u^R - w)\|_{L^2(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})}^2 \\
& = \min_{v \in \mathcal{P}} \left(\left\| v + \frac{1}{2} \sum_{k=0}^p \frac{J_i^k}{k!} (x - x_i)^k \right\|_{L^2(\Delta_i^-)}^2 + \left\| v - \frac{1}{2} \sum_{k=0}^p \frac{J_i^k}{k!} (x - x_i)^k \right\|_{L^2(\Delta_i^+)}^2 \right) \\
& = h^{2p+3} \min_{\hat{v} \in \mathcal{P}} \left(\left\| \hat{v} + \frac{1}{2} \sum_{k=0}^p \frac{D_i^k}{k!} \xi^k \right\|_{L^2(-\frac{1}{2}, 0)}^2 + \left\| \hat{v} - \frac{1}{2} \sum_{k=0}^p \frac{D_i^k}{k!} \xi^k \right\|_{L^2(0, \frac{1}{2})}^2 \right) \\
& = h^{2p+3} Q(D_i^0, D_i^1, \dots, D_i^p). \tag{2.4}
\end{aligned}$$

Plugging (2.4) into (2.3), we have proven (2.1). #

Theorem 2.4 *If $u \in W_1^{p+1}(a, b)$, and u^R is as described previously, then, there are positive constants C_1 and C_{12}^p , independent of h , u and u^R , such that*

$$\|u - u^R\|_{L^1(a, b)} \geq h^{p+1} \left[C_{12}^p \sum_{1 \leq i \leq N-1} h \sqrt{Q(D_i^0, D_i^1, \dots, D_i^p)} - C_1 |u|_{W_1^{p+1}(a, b)} \right]. \tag{2.5}$$

Proof. Since $u \in W^{p+1,1}(a, b)$, there is a piecewise polynomial $u^I \in \mathcal{P}_c$, such that

$$\|u - u^I\|_{L^1(a, b)} \leq C_1 h^{p+1} |u|_{W_1^{p+1}(a, b)},$$

where $C_1 > 0$ is a constant independent of h and u . Obviously,

$$\begin{aligned}
& \|u - u^R\|_{L^1(a, b)} \geq \|u^I - u^R\|_{L^1(a, b)} - \|u - u^I\|_{L^1(a, b)} \\
& \geq \sum_{1 \leq i \leq N-1} \|u^I - u^R\|_{L^1(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})} - C_1 h^{p+1} |u|_{W_1^{p+1}(a, b)}. \tag{2.6}
\end{aligned}$$

By the standard scaling argument, and (2.4) which remains valid for the current u^I , it is easy to prove that there is a constant C_{12}^p such that

$$\|u^I - u^R\|_{L^1(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})} \geq C_{12}^p h^{\frac{1}{2}} \|u^I - u^R\|_{L^2(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})} \geq C_{12}^p h^{p+2} \sqrt{Q(D_i^0, D_i^1, \dots, D_i^p)}.$$

Plugging this into (2.6), we have (2.5) proven. #

Theorem 2.5 *If $u \in W_\infty^{p+1}[a, b]$, and u^R is as described previously, then, there is a constant $C_\infty > 0$, independent of h , u and u^R , such that*

$$\|u - u^R\|_{L^\infty(a, b)} \geq h^{p+1} \left[\max_{1 \leq i \leq N-1} \sqrt{Q(D_i^0, D_i^1, \dots, D_i^p)} - C_\infty |u|_{W_\infty^{p+1}[a, b]} \right]. \tag{2.7}$$

Proof. Since $u \in W_\infty^{p+1}[a, b]$, there is a piecewise polynomial $u^I \in \mathcal{P}_c$, such that

$$\|u - u^I\|_{L^\infty(a, b)} \leq C_\infty h^{p+1} |u|_{W_\infty^{p+1}[a, b]},$$

where $C_\infty > 0$ is a constant independent of h and u . Obviously,

$$\begin{aligned} \|u - u^R\|_{L^\infty(a, b)} &\geq \|u^I - u^R\|_{L^\infty(a, b)} - \|u - u^I\|_{L^\infty(a, b)} \\ &\geq \max_{1 \leq i \leq N-1} \|u^I - u^R\|_{L^\infty(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})} - C_\infty h^{p+1} |u|_{W_\infty^{p+1}[a, b]}. \end{aligned} \quad (2.8)$$

Now, let $U_i = \|u^I - u^R\|_{L^\infty(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})}$, then, by using (2.4) for the current u^I ,

$$\begin{aligned} U_i &= \sqrt{\int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \frac{1}{h} U_i^2 dx} \geq \sqrt{\int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \frac{1}{h} (u^I - u^R)^2 dx} \\ &= \sqrt{\frac{1}{h} \|u^I - u^R\|_{L^2(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})}^2} \geq h^{p+1} \sqrt{Q(D_i^0, D_i^1, \dots, D_i^p)}. \end{aligned} \quad (2.9)$$

By plugging (2.9) into (2.8), we have proven (2.7). #

3 Conclusion remarks

According to Theorem 2.3, in order to have

$$\|u - u^R\|_{L^2(a, b)} \leq \mathcal{O}(h^{p+1})$$

for any function $u \in H^{p+1}(a, b)$, u^R must be numerically H^{p+1} -smooth. That is

$$\sum_{i=1}^{N-1} h Q(D_i^0, D_i^1, \dots, D_i^p) \leq \mathcal{O}(1). \quad (3.1)$$

According to Theorem 2.4, in order to have

$$\|u - u^R\|_{L^1(a, b)} \leq \mathcal{O}(h^{p+1})$$

for any function $u \in W_1^{p+1}(a, b)$, u^R must be numerically W_1^{p+1} -smooth. That is

$$\sum_{i=1}^{N-1} h \sqrt{Q(D_i^0, D_i^1, \dots, D_i^p)} \leq \mathcal{O}(1). \quad (3.2)$$

According to Theorem 2.5, in order to have

$$\|u - u^R\|_{L^\infty(a, b)} \leq \mathcal{O}(h^{p+1})$$

for any function $u \in W_\infty^{p+1}[a, b]$, u^R must be numerically W_∞^{p+1} -smooth. That is

$$Q(D_i^0, D_i^1, \dots, D_i^p) \leq \mathcal{O}(1) \quad (3.3)$$

for all i . Because $Q(d_0, d_1, \dots, d_p)$ is a positive definite quadratic form, the last inequality is equivalent to that $|D_i^k| \leq \mathcal{O}(1)$ for all k and i . This is what we mean by the necessity of the numerical smoothness. To be more explicit, $|D_i^k| \leq \mathcal{O}(1)$ is equivalent to $|J_i^k| \leq \mathcal{O}(h^{p+1-k})$. That is, the jumps of the k -th derivative of u^R need to be as small as $\mathcal{O}(h^{p+1-k})$, for $k = 0, 1, \dots, p$.

Although the results and the proofs of the theorems seem to be very simple, the author believes that the impact of the theorems could be far reaching. For a numerical solution of a time-dependent PDE, the theorem implies that a fully-discrete scheme must enforce numerical smoothness, otherwise the convergence of optimal rate must have been lost. Consequently, some of the traditional numerical stability notions are possibly inadequate, if they do not enforce numerical smoothness.

References

- [1] S. C. Brenner and L. R. Scott, *The Mathematical Theory of Finite Element Methods*, 3rd Edition, Texts in Applied Mathematics, Springer, New York, 2008.
- [2] T. Sun and D. Rumsey, *Numerical smoothness and error analysis for RKDG on the scalar nonlinear conservation laws*, originally submitted for publication in September, 2010; arXiv:1105.1393, May 2011.
- [3] T. Sun, *Numerical smoothness and error analysis for WENO on the scalar nonlinear conservation laws*, submitted for publication in NMPDE, October, 2011.
- [4] Q. Zhang and C.-W. Shu, *Stability analysis and a priori error estimates to the third order explicit Runge-Kutta discontinuous Galerkin Method for scalar conservation laws*, SIAM Journal on Numerical Analysis, 48 (2010), 1038-1063.